
Gaussian Processes for Classification & Brief Summary of the Course

*Prof. Nicholas Zabaras
University of Notre Dame
Notre Dame, IN, USA*

*Email: nzabaras@gmail.com
URL: <https://www.zabaras.com/>*

November 15, 2017



Contents

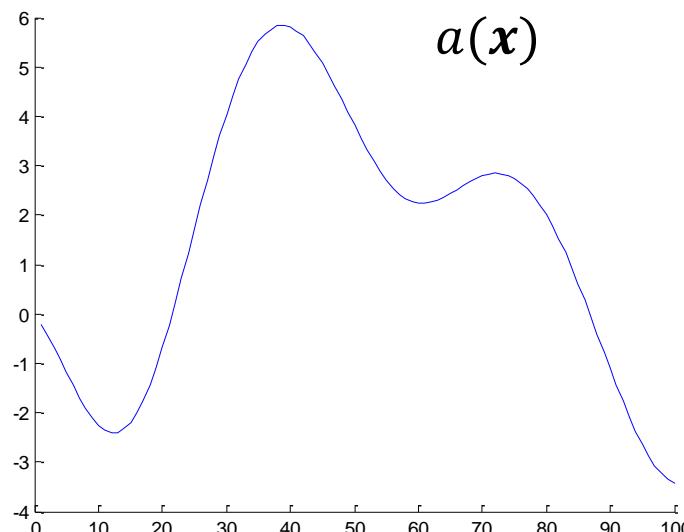
- Gaussian Process Classification, Laplace Approximation, Connections to Bayesian Neural Networks
- Markov and Chebyshev Inequalities, Law of Large Numbers
- MLE for Gaussians, Robbins-Monro, Curse of Dimensionality
- Bayes' theorem and Gaussian Linear Models, Posterior Inference and Prediction, Marginal Likelihood
- Exponential family of distributions, Conjugate priors
- Empirical Bayes, Evidence approximation, Model Validation and Bayes factors, Occam's Razor
- Laplace approximation
- Sampling from an arbitrary distribution, Inverse Method, Rejection Sampling, Importance Sampling, Gibbs Sampling, MCMC, Metropolis-Hastings, Sequential Importance Sampling and Particle Methods, Reversible Jump MCMC
- Latent Variables and Expectation-Maximization
- Probabilistic PCA and Generative View

Following closely: Bishop CM, *Pattern Recognition and Machine Learning*, Springer, 2006 (Chapter 6)

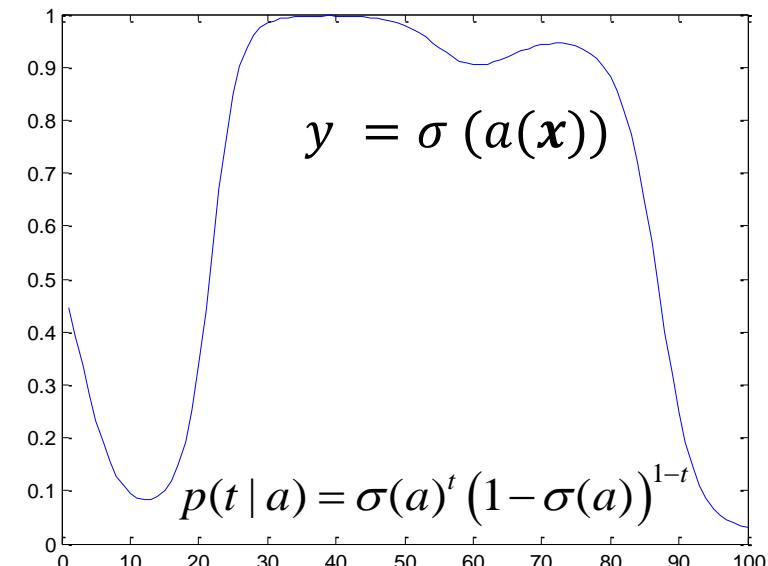


Gaussian Processes for Classification

- Objective: model posterior probabilities of the target variable for a new input
- Problem: we need to map values to interval $(0, 1)$
- Solution: use a Gaussian process together with a non-linear activation function (e.g., sigmoid)



Sample from a Gaussian process over $a(x)$

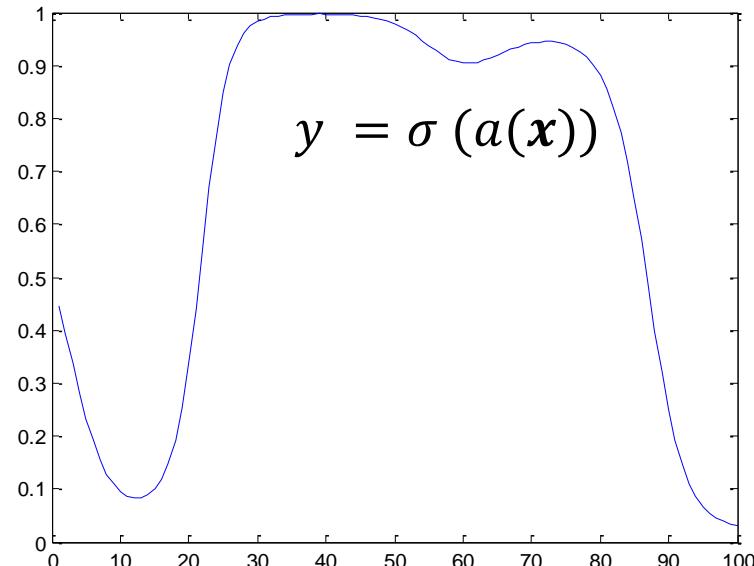
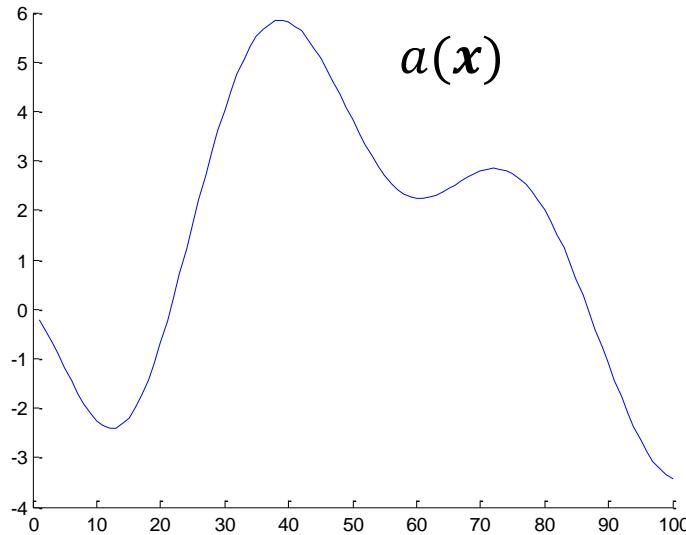


MatLab Code



Gaussian Process for Classification

- Consider a two-class problem with target values $t \in \{0, 1\}$. Define a Gaussian process over a function $a(x)$ and transform a using the logistic sigmoid to $y = \sigma(a(x))$.
- This will give a non-Gaussian stochastic process over functions $y(x)$ where $y \in \{0, 1\}$.



[MatLab Code](#)

- The probability over the target t is a Bernoulli: $p(t | a) = \sigma(a)^t (1 - \sigma(a))^{1-t}$

Gaussian Process for Classification

- We denote the training set inputs by $\mathbf{x}_1, \dots, \mathbf{x}_N$ with corresponding observed target variables $\mathbf{t} = (t_1, \dots, t_N)^T$.
- We also consider a single test point \mathbf{x}_{N+1} with target value t_{N+1} .
- Our goal is to determine the predictive distribution $p(t_{N+1}|\mathbf{t})$, where we left the conditioning on the input variables implicit.
- Introduce a Gaussian process prior over the vector \mathbf{a}_{N+1} , which has components $a(\mathbf{x}_1), \dots, a(\mathbf{x}_{N+1})$.
- This defines a non-Gaussian process over \mathbf{t}_{N+1} , and by conditioning on \mathbf{t}_N we obtain the required predictive distribution $p(t_{N+1}|\mathbf{t})$.
- The Gaussian process prior for \mathbf{a}_{N+1} takes the form

$$p(\mathbf{a}_{N+1}) = \mathcal{N}(\mathbf{a}_{N+1} | \mathbf{0}, \mathbf{C}_{N+1})$$



Gaussian Process for Classification

- The covariance matrix does not include a noise term because we assume that all training data points are correctly labeled.
- However, for numerical reasons it is convenient to introduce a noise-like term governed by a parameter ν that ensures that the covariance matrix is positive definite. Thus the covariance matrix C_{N+1} has elements given by

$$C(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \nu \delta_{nm}$$

- The kernel k is positive semidefinite and ν is fixed. $k(\mathbf{x}, \mathbf{x}')$ is governed by parameters θ that we will compute from the training data.
- For two class problems, it is sufficient to predict:

$$p(t_{N+1} = 1 | \mathbf{t}_N)$$

Gaussian Process for Classification

- The predictive distribution is given as:

$$\begin{aligned} p(t_{N+1} = 1 | \mathbf{t}_N) &= \int p(t_{N+1} = 1, a_{N+1} | \mathbf{t}_N) da_{N+1} \\ &= \int p(t_{N+1} = 1 | a_{N+1}) p(a_{N+1} | \mathbf{t}_N) da_{N+1} = \int \sigma(a_{N+1}) p(a_{N+1} | \mathbf{t}_N) da_{N+1} \end{aligned}$$

- The integral is analytically intractable. Can be approximated with sampling (Neal)
- Integration can be analytical if $p(a_{N+1} | t_N)$ can be approximated with a Gaussian (Williams & Barber, Gibbs & MacKay, Gibbs)
- Approximations based on EP can also be used (Opper & Winther, Minka, Seeger)
 - Neal, R. M. (1997). [Monte Carlo implementation of Gaussian process models for Bayesian regression and classification](#). [Technical Report 9702](#), Department of Computer Statistics, University of Toronto.
 - Williams, C. K. I. and D. Barber (1998). [Bayesian classification with Gaussian processes](#). [IEEE Transactions on Pattern Analysis and Machine Intelligence](#) **20**, 1342–1351.
 - Gibbs, M. N. and D. J. C. MacKay (2000). [Variational Gaussian process classifiers](#). [IEEE Transactions on Neural Networks](#) **11**, 1458–1464.
 - Gibbs, M. N. (1997). [Bayesian Gaussian processes for regression and classification](#). Phd thesis, University of Cambridge.
 - Opper, M. and O. Winther (2000b). [Gaussian processes for classification](#). [Neural Computation](#) **12**(11), 2655–2684.
 - Minka, T. (2001b). [A family of approximate algorithms for Bayesian inference](#). Ph. D. thesis, MIT.
 - Seeger, M. (2003). [Bayesian Gaussian Process Models: PAC-Bayesian Generalization Error Bounds and Sparse Approximations](#). Ph. D. thesis, University of Edinburg.

Gaussian Process for Classification

- We need to compute $p(t_{N+1} = 1 | \mathbf{t}_N) = \int \sigma(a_{N+1}) p(a_{N+1} | \mathbf{t}_N) da_{N+1}$
- One can compute the convolution of a Gaussian and a sigmoid function:

$$\int \sigma(a) \mathcal{N}(a|\mu, \sigma^2) da \simeq \sigma(\kappa(\sigma^2)\mu)$$

$$\kappa(\sigma^2) = \sqrt{1 + \frac{\pi\sigma^2}{8}}$$

We can thus try to approximate the posterior distribution $p(a_{N+1} | \mathbf{t}_N)$ as Gaussian

- How do we approximate the posterior?
 - Variational inference and make use of the local variational bound on the logistic sigmoid.
 - Expectation propagation
 - Laplace approximation



Laplace Approximation

- We can rewrite the posterior over a_{N+1} using Bayes' theorem:

$$\begin{aligned} p(a_{N+1} | \mathbf{t}_N) &= \int p(a_{N+1}, \mathbf{a}_N | \mathbf{t}_N) d\mathbf{a}_N \\ &= \frac{1}{p(\mathbf{t}_N)} \int p(a_{N+1}, \mathbf{a}_N) p(\mathbf{t}_N | a_{N+1}, \mathbf{a}_N) d\mathbf{a}_N \\ &= \frac{1}{p(\mathbf{t}_N)} \int p(a_{N+1} / \mathbf{a}_N) p(\mathbf{a}_N) p(\mathbf{t}_N | \mathbf{a}_N) d\mathbf{a}_N \\ &= \int p(a_{N+1} / \mathbf{a}_N) p(\mathbf{a}_N | \mathbf{t}_N) d\mathbf{a}_N \end{aligned}$$

- We know how to compute mean and covariance for $p(\mathbf{a}_N | \mathbf{a}_N)$.

$$p(a_{N+1} | \mathbf{a}_N) = \mathcal{N}\left(\mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{a}_N, c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}\right)$$

- It remains to find a Gaussian approximation only for $p(\mathbf{a}_N | \mathbf{t}_N)$
- This is done noting that $p(\mathbf{a}_N | \mathbf{t}_N) \propto p(\mathbf{t}_N | \mathbf{a}_N) p(\mathbf{a}_N)$

Laplace Approximation

- We rewrote the posterior over a_{N+1} using Bayes' theorem:

$$p(a_{N+1} | \mathbf{t}_N) = \int p(a_{N+1} | \mathbf{a}_N) p(\mathbf{a}_N | \mathbf{t}_N) d\mathbf{a}_N$$

where

$$p(a_{N+1} | \mathbf{a}_N) = \mathcal{N}\left(\mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{a}_N, c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}\right)$$

- Also using $p(\mathbf{a}_N | \mathbf{t}_N) \propto p(\mathbf{t}_N | \mathbf{a}_N) p(\mathbf{a}_N)$, $p(\mathbf{a}_N) = \mathcal{N}(0, \mathbf{C}_N)$ and

$$p(\mathbf{t}_N | \mathbf{a}_N) = \prod_{n=1}^N \sigma(a_n)^{t_n} (1 - \sigma(a_n))^{1-t_n} = \prod_{n=1}^N e^{a_n t_n} \sigma(-a_n)$$

where in the last derivation we used: $\sigma(a_n) = \frac{1}{1 + e^{-a_n}}$ and

$$\sigma(a_n)^{t_n} (1 - \sigma(a_n))^{1-t_n} = \underbrace{(1 - \sigma(a_n))}_{\sigma(-a_n)} \left(\underbrace{\frac{\sigma(a_n)}{1 - \sigma(a_n)}}_{e^{a_n}} \right)^{t_n} = e^{a_n t_n} \sigma(-a_n)$$

Laplace Approximation

- The log of $p(\mathbf{a}_N / \mathbf{t}_N) \propto p(\mathbf{t}_N / \mathbf{a}_N)p(\mathbf{a}_N)$ which up to an additive constant is:

$$\begin{aligned} p(\mathbf{a}_N / \mathbf{t}_N) &\propto p(\mathbf{a}_N)p(\mathbf{t}_N / \mathbf{a}_N) \\ &\propto \mathcal{N}(0, \mathbf{C}_N) \prod_{n=1}^N e^{a_n t_n} \underbrace{\sigma(-a_n)}_{1/(1+e^{a_n})} \Rightarrow \end{aligned}$$

$$\begin{aligned} \Psi(\mathbf{a}_N) &= \ln p(\mathbf{a}_N / \mathbf{t}_N) = \ln p(\mathbf{t}_N / \mathbf{a}_N) + \ln p(\mathbf{a}_N) \\ &= -\frac{1}{2} \mathbf{a}_N^T \mathbf{C}_N^{-1} \mathbf{a}_N - \frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{C}_N| + \mathbf{t}_N^T \mathbf{a}_N - \sum_{n=1}^N \ln(1 + e^{a_n}) \end{aligned}$$

- We will need to compute the 2nd derivative of this to come up with the Laplace approximation.

Laplace Approximation

$$\Psi(\mathbf{a}_N) = -\frac{1}{2}\mathbf{a}_N^T \mathbf{C}_N^{-1} \mathbf{a}_N - \frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{C}_N| + \mathbf{t}_N^T \mathbf{a}_N - \sum_{n=1}^N \ln(1 + e^{a_n})$$

- The gradients are given as

$$\nabla \Psi(\mathbf{a}_N) = \mathbf{t}_N - \boldsymbol{\sigma}_N - \mathbf{C}_N^{-1} \mathbf{a}_N, \quad \boldsymbol{\sigma}_N = (\sigma(a_1), \dots, \sigma(a_N))^T$$

$$\nabla^2 \Psi(\mathbf{a}_N) = -\mathbf{W}_N - \mathbf{C}_N^{-1}, \quad \mathbf{W}_N = \text{diag}(\sigma(a_1)(1 - \sigma(a_1)), \dots, \sigma(a_N)(1 - \sigma(a_N)))$$

- We cannot find the mode by setting $\nabla \Psi(\mathbf{a}_N) = 0$ since $\boldsymbol{\sigma}_N$ depends nonlinearly on \mathbf{a}_N . For this we will use the **IRLS algorithm**.
- Note that $0 < \sigma(a_1)(1 - \sigma(a_1)) < 1/4 \Rightarrow \mathbf{W}_N$ is positive definite.
- The Hessian $A = -\nabla^2 \Psi(\mathbf{a}_N)$ is positive definite and thus the posterior $p(\mathbf{a}_N / \mathbf{t}_N)$ is log convex and has a single global mode.
- At the mode, the gradient vanishes giving: $\mathbf{a}_N^* = \mathbf{C}_N (\mathbf{t}_N - \boldsymbol{\sigma}_N)$

Iterative Reweighted Least Squares

- We can use the iterative reweighted least squares (IRLS) algorithm ($\mathbf{w}^{new} = \mathbf{w}^{old} - \mathbf{H}^{-1} \nabla E(\mathbf{w})$) to update for \mathbf{a}_N (find the minimum of $-\ln p(\mathbf{a}_N / \mathbf{t}_N)$).
- The iterative update equation for \mathbf{a}_N is:

$$\mathbf{a}_N^{new} = \mathbf{C}_N (\mathbf{I} + \mathbf{W}_N \mathbf{C}_N)^{-1} (\mathbf{t}_N - \boldsymbol{\sigma}_n + \mathbf{W}_N \mathbf{a}_N)$$

Here we used

$$\begin{aligned}\mathbf{a}_N^{new} &= \mathbf{a}_N + (\mathbf{W}_N + \mathbf{C}_N^{-1})^{-1} (\mathbf{t}_N - \boldsymbol{\sigma}_N - \mathbf{C}_N^{-1} \mathbf{a}_N) = \mathbf{a}_N + \mathbf{C}_N (\mathbf{W}_N \mathbf{C}_N + \mathbf{I})^{-1} (\mathbf{t}_N - \boldsymbol{\sigma}_N - \mathbf{C}_N^{-1} \mathbf{a}_N) \\ &= \mathbf{C}_N (\mathbf{W}_N \mathbf{C}_N + \mathbf{I})^{-1} \{ (\mathbf{W}_N \mathbf{C}_N + \mathbf{I}) \mathbf{C}_N^{-1} \mathbf{a}_N + (\mathbf{t}_N - \boldsymbol{\sigma}_N - \mathbf{C}_N^{-1} \mathbf{a}_N) \} \\ &= \mathbf{C}_N (\mathbf{W}_N \mathbf{C}_N + \mathbf{I})^{-1} \{ \mathbf{W}_N \mathbf{a}_N + \mathbf{t}_N - \boldsymbol{\sigma}_N \}\end{aligned}$$

Laplace Approximation

- We can use the iterative reweighted least squares (IRLS) algorithm ($\mathbf{w}^{new} = \mathbf{w}^{old} - \mathbf{H}^{-1}\nabla E(\mathbf{w})$) to update for \mathbf{a}_N .
- The iterative update equation for \mathbf{a}_N is:

$$\mathbf{a}_N^{new} = \mathbf{C}_N (\mathbf{I} + \mathbf{W}_N \mathbf{C}_N)^{-1} (\mathbf{t}_n - \boldsymbol{\sigma}_n + \mathbf{W}_N \mathbf{a}_N)$$

- The mode position \mathbf{a}_N^* and the Hessian matrix \mathbf{H} at this position $\mathbf{H} = -\nabla^2 \Psi(\mathbf{a}_N) = \mathbf{W}_N + \mathbf{C}_N^{-1}$ define the Gaussian approximation to $p(\mathbf{a}_N | \mathbf{t}_n)$

$$q(\mathbf{a}_N) = \mathcal{N}(\mathbf{a}_N | \mathbf{a}_N^*, \mathbf{H}^{-1})$$

In the Hessian note that \mathbf{W}_N is computed using \mathbf{a}_N^* .



Laplace Approximation

- Now we can go back to the formulas for linear Gaussian models and compute the integrals

$$p(a_{N+1} | \mathbf{t}_N) = \int p(a_{N+1} | \mathbf{a}_N) p(\mathbf{a}_N | \mathbf{t}_N) d\mathbf{a}_N$$

$$p(a_{N+1} | \mathbf{a}_N) = \mathcal{N}\left(\mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{a}_N, c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}\right)$$

$$p(\mathbf{a}_N | \mathbf{t}_N) \approx q(\mathbf{a}_N) = \mathcal{N}(\mathbf{a}_N | \mathbf{a}_N^*, \mathbf{H}^{-1}), \mathbf{H} = \mathbf{W}_N + \mathbf{C}_N^{-1}$$

from which we have:

$$p(a_{N+1} | \mathbf{t}_N) = \mathcal{N}\left(\mathbf{k}^T (\mathbf{t}_N - \boldsymbol{\sigma}_N), c - \mathbf{k}^T (\mathbf{W}_N^{-1} + \mathbf{C}_N)^{-1} \mathbf{k}\right)$$

- Finally, we can also compute $p(t_{N+1} | \mathbf{t})$ (or the decision boundary $p(t_{N+1} | \mathbf{t}) = 0.5$ – in this case $\mu = \mathbf{k}^T (\mathbf{t}_N - \boldsymbol{\sigma}_N) = 0$)

$$p(t_{N+1} = 1 | \mathbf{t}_N) = \int \sigma(a_{N+1}) p(a_{N+1} | \mathbf{t}_N) da_{N+1}$$

using $\int \sigma(a) \mathcal{N}(a | \mu, \sigma^2) da \approx \sigma(\kappa(\sigma^2)\mu)$, $\kappa(\sigma^2) = \sqrt{1 + \frac{\pi\sigma^2}{8}}$

Learning the Hyperparameters

- To determine the parameters θ of the covariance function, we can maximize the likelihood function:

$$p(\mathbf{t}_N | \theta) = \int p(\mathbf{t}_N | \mathbf{a}_N) p(\mathbf{a}_N | \theta) d\mathbf{a}_N = \int p(\mathbf{t}_N, \mathbf{a}_N | \theta) d\mathbf{a}_N$$

- The integral is analytically intractable. The Laplace approximation can be applied again:
$$\ln p(\mathbf{t}_N, \mathbf{a}_N) = \Psi(\mathbf{a}_N) = -\frac{1}{2} \mathbf{a}_N^T \mathbf{C}_N^{-1} \mathbf{a}_N - \frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{C}_N| + \mathbf{t}_N^T \mathbf{a}_N - \sum_{n=1}^N \ln(1 + e^{a_n})$$
$$\ln p(\mathbf{t}_N, \mathbf{a}_N) \approx \Psi(\mathbf{a}_N^*) - \frac{1}{2} (\mathbf{a}_N - \mathbf{a}_N^*)^T (\mathbf{W}_N + \mathbf{C}_N^{-1}) (\mathbf{a}_N - \mathbf{a}_N^*)$$
- Integrating in \mathbf{a}_N and using the normalization factor of a Gaussian gives:

$$\ln p(\mathbf{t}_N | \theta) = \Psi(\mathbf{a}_N^*) - \frac{1}{2} \ln |\mathbf{W}_N + \mathbf{C}_N^{-1}| + \frac{N}{2} \ln(2\pi)$$

where

$$\Psi(\mathbf{a}_N^*) = \ln p(\mathbf{a}_N^* | \theta) + \ln p(\mathbf{t}_N | \mathbf{a}_N^*)$$

- We now need an expression for the gradient of the log of $p(\mathbf{t}_N | \theta)$. Both \mathbf{a}_N^* and \mathbf{C}_N depend on θ

Learning the Hyperparameters

- To determine the parameters , maximize the likelihood:

$$\Psi(\mathbf{a}_N^*) = -\frac{1}{2} \mathbf{a}_N^{*T} \mathbf{C}_N^{-1} \mathbf{a}_N^* - \frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{C}_N| + \mathbf{t}_N^T \mathbf{a}_N^* - \sum_{n=1}^N \ln(1 + e^{a_n^*})$$

$$\ln p(\mathbf{t}_N | \boldsymbol{\theta}) = \Psi(\mathbf{a}_N^*) - \frac{1}{2} \ln |\mathbf{W}_N + \mathbf{C}_N^{-1}| + \frac{N}{2} \ln(2\pi)$$

- Using $\frac{\partial}{\partial x} \ln |\mathbf{A}| = \text{Tr}\left(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x}\right)$ and $\frac{\partial \mathbf{A}^{-1}}{\partial x} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1}$, we obtain

the following terms of $\frac{\partial}{\partial \theta_j} \ln p(\mathbf{t}_N | \boldsymbol{\theta})$ with explicit dependence on $\boldsymbol{\theta}$:

$$\frac{1}{2} \mathbf{a}_N^{*T} \mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_j} \mathbf{C}_N^{-1} \mathbf{a}_N^* - \frac{1}{2} \text{Tr} \left[\mathbf{C}_N^{-1} \left[\mathbf{I} - (\mathbf{W}_N + \mathbf{C}_N^{-1})^{-1} \mathbf{C}_N^{-1} \right] \frac{\partial \mathbf{C}_N}{\partial \theta_j} \right] =$$

$$\frac{1}{2} \mathbf{a}_N^{*T} \mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_j} \mathbf{C}_N^{-1} \mathbf{a}_N^* - \frac{1}{2} \text{Tr} \left[\mathbf{C}_N^{-1} \left[(\mathbf{C}_N \mathbf{W}_N + \mathbf{I})(\mathbf{C}_N \mathbf{W}_N + \mathbf{I})^{-1} - (\mathbf{C}_N \mathbf{W}_N + \mathbf{I})^{-1} \right] \frac{\partial \mathbf{C}_N}{\partial \theta_j} \right]$$

$$\frac{1}{2} \mathbf{a}_N^{*T} \mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_j} \mathbf{C}_N^{-1} \mathbf{a}_N^* - \frac{1}{2} \text{Tr} \left[(\mathbf{I} + \mathbf{C}_N \mathbf{W}_N)^{-1} \mathbf{W}_N \frac{\partial \mathbf{C}_N}{\partial \theta_j} \right]$$



Learning the Hyperparameters

- To determine the parameters, maximize the likelihood:

$$\Psi(\mathbf{a}_N^*) = -\frac{1}{2} \mathbf{a}_N^{*T} \mathbf{C}_N^{-1} \mathbf{a}_N^* - \frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{C}_N| + \mathbf{t}_N^T \mathbf{a}_N^* - \sum_{n=1}^N \ln(1 + e^{a_n^*})$$

$$\ln p(\mathbf{t}_N | \boldsymbol{\theta}) = \Psi(\mathbf{a}_N^*) - \frac{1}{2} \ln |\mathbf{W}_N + \mathbf{C}_N^{-1}| + \frac{N}{2} \ln(2\pi)$$

- Using $\nabla \Psi(\mathbf{a}_N^*) = 0$, the term in $\frac{\partial}{\partial \theta_j} \ln p(\mathbf{t}_N | \boldsymbol{\theta})$ with explicit dependence on a_n^* is as follows:

$$-\frac{1}{2} \sum_{n=1}^N \frac{\partial \ln |\mathbf{W}_N + \mathbf{C}_N^{-1}|}{\partial a_n^*} \frac{\partial a_n^*}{\partial \theta_j}$$

- Using $\frac{\partial}{\partial x} \ln |\mathbf{A}| = \text{Tr}\left(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x}\right)$ and $\mathbf{W}_N = \text{diag}(\sigma(a_i)(1 - \sigma(a_i)))$, it becomes $d\sigma / da_i = \sigma(a_i)(1 - \sigma(a_i))$

$$-\frac{1}{2} \sum_{n=1}^N \text{Tr} \left[(\mathbf{W}_N + \mathbf{C}_N^{-1})^{-1} \underbrace{\frac{\partial \mathbf{W}_N}{\partial \sigma(a_n)}}_{(1-2\sigma_n^*)} \underbrace{\frac{\partial \sigma(a_n)}{\partial a_n}}_{\sigma_n^*(1-\sigma_n^*)} \right] \frac{\partial a_n^*}{\partial \theta_j} = -\frac{1}{2} \sum_{n=1}^N \left[(\mathbf{I} + \mathbf{C}_N \mathbf{W}_N)^{-1} \mathbf{C}_N \right]_{nn} \sigma_n^* (1 - \sigma_n^*) (1 - 2\sigma_n^*) \frac{\partial a_n^*}{\partial \theta_j}$$

Learning Hyperparameters

- We finally can compute $\partial \mathbf{a}_n^* / \partial \theta_j$ as follows : $\mathbf{a}_N^* = \mathbf{C}_N (\mathbf{t}_N - \boldsymbol{\sigma}_N) \Rightarrow$

$$\frac{\partial \mathbf{a}_N^*}{\partial \theta_j} = \frac{\partial \mathbf{C}_N}{\partial \theta_j} (\mathbf{t}_N - \boldsymbol{\sigma}_N) - \mathbf{C}_N \mathbf{W}_N \frac{\partial \mathbf{a}_N^*}{\partial \theta_j}$$

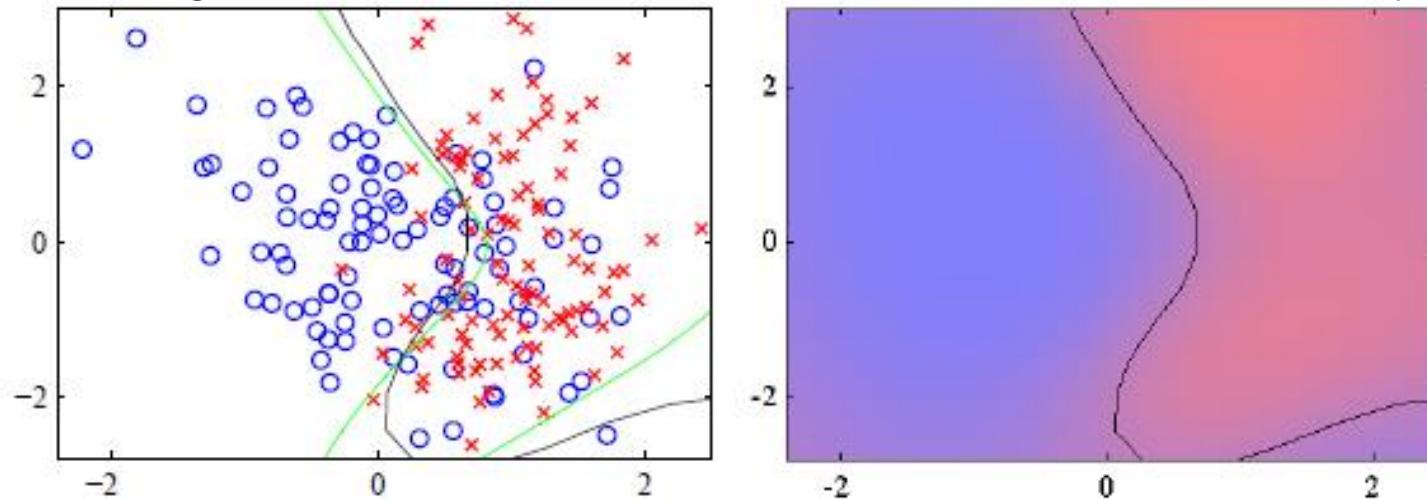
$\mathbf{W}_N = \text{diag}(\sigma(a_i)(1-\sigma(a_i)))$
 $d\sigma / da_i = \sigma(a_i)(1-\sigma(a_i))$

$$\frac{\partial \mathbf{a}_N^*}{\partial \theta_j} = (\mathbf{I} + \mathbf{C}_N \mathbf{W}_N)^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_j} (\mathbf{t}_N - \boldsymbol{\sigma}_N)$$

- We now have all the terms in the gradient $\frac{\partial}{\partial \theta_j} \ln p(\mathbf{t}_N | \theta)$ and can employ gradient optimization techniques.

Binary Classifier Based on Gaussian Process

- We now have a binary classifier based on a Gaussian process. On left the data together with the optimal decision boundary from the true distribution in green, and the decision boundary from the Gaussian process classifier in black.
- On right the predicted posterior probability for the blue and red classes together with the Gaussian process decision boundary.



Requires installing [NetLab](#)
and [setting the appropriate path](#)

[MatLab Code](#)

- Williams, C. K. I. and D. Barber (1998). [Bayesian classification with Gaussian processes](#). *IEEE Transactions on Pattern Analysis and Machine Intelligence* **20**, 1342–1351.



Connections to Neural Networks

□ Neural Networks

- The range of representable functions is governed by the number M of hidden units
- Within the maximum likelihood framework, they overfit as M comes close to the number of training samples

□ Bayesian Neural Networks

- The prior over w in conjunction with the network function $f(x, w)$ produces a prior distribution over functions from $y(x)$
- The distribution of functions will tend to a GP as $M \rightarrow \infty$.
 - In this limit the output variables of the neural network become independent.
 - However, the property that the outputs share hidden units and thus all of them affecting the values of the weights is lost in this limit.

▪ Neal, R. M. (1996). *Bayesian Learning for Neural Networks*. Springer. Lecture Notes in Statistics 118.



Connections to Neural Networks

- There are explicit forms for the covariance in the case of probit and Gaussian hidden unit activation functions.
 - These kernel functions $k(x, x')$ are nonstationary, i.e. not a function of $x - x'$.
 - This is a consequence of the Gaussian weight prior being centered on zero which breaks translation invariance in weight space.
- By working directly with the covariance function we have implicitly marginalized over the distribution of weights.
- If the weight prior is governed by hyperparameters, then their values will determine the length scales of the distribution over functions.

▪ [Williams, C. K.](#) I. (1998). [Computation with infinite neural networks](#). [Neural Computation 10\(5\), 1203–1216](#).

Statistical Computing, University of Notre Dame, Notre Dame, IN, USA (Fall 2017, N. Zabaras)



Course Summary



Markov and Chebyshev Inequalities

- You can show ([Markov's inequality](#)) that if X is a non-negative integrable random variable and for any $a > 0$:

$$\Pr[X \geq a] \leq \frac{\mathbb{E}[X]}{a}$$

Indeed : $\mathbb{E}[X] \geq \int_a^\infty x\pi(x)dx \geq a \int_a^\infty \pi(x)dx = a \Pr[X \geq a]$

- You can generalize this using any function of the random variable X as:

$$\Pr[f(X) \geq a] \leq \frac{\mathbb{E}[f(X)]}{a}$$

- Using $f(X) = (X - \mathbb{E}[X])^2$, we derive the following Chebyshev inequality:

$$a \equiv \varepsilon^2, \mathbb{E}[(X - \mathbb{E}[X])^2] = \sigma^2$$

$$\Pr[|X - \mathbb{E}[X]| \geq \varepsilon] \leq \frac{\sigma^2}{\varepsilon^2} \quad \forall \varepsilon$$

- In terms of std's, we can restate as : $\Pr[|X - \mathbb{E}[X]| \geq \varepsilon\sigma] \leq \frac{1}{\varepsilon^2} \quad \forall \varepsilon$

- Thus the probability of X being more than 2σ away from $\mathbb{E}[X]$ is $\leq \frac{1}{4}$.



The Law of Large Numbers (LLN)

- Let X_i for $i = 1, 2, \dots, n$ be independent and identically distributed random variables (i.i.d.) with $\mathbb{E}(X_i) = \mu$ and variance $\text{Var}(X_i) = \sigma^2$.

- Let $\bar{X}_n = \sum_{i=1}^n \frac{X_i}{n}$

- Note that $\mathbb{E}[\bar{X}_n] = \sum_{i=1}^n \frac{\mu}{n} = \mu$

$$\text{Var}[\bar{X}_n] = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{\sigma^2}{n}$$

- Weak LLN: $\lim_{n \rightarrow \infty} \Pr[|\bar{X}_n - \mu| \geq \varepsilon] = 0 \quad \forall \varepsilon > 0$

- Strong LLN: $\lim_{n \rightarrow \infty} \bar{X}_n = \mu$ almost surely

This means that with probability one, the average of any realizations of x_1, x_2, \dots of the random variables X_1, X_2, \dots converges to the mean.



MLE for a Multivariate Gaussian

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n \equiv \bar{x}, \quad \Sigma_{ML} = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})(x_n - \mu_{ML})^T = \frac{1}{N} \sum_{n=1}^N x_n x_n^T - \bar{x} \bar{x}^T$$

- Note that *the unconstrained maximization of the log-likelihood gives a symmetric Σ .*
- As for the univariate case, we can define an **unbiased covariance** as:

$$\bar{\Sigma}_{ML} = \frac{1}{N-1} \sum_{n=1}^N (x_n - \mu_{ML})(x_n - \mu_{ML})^T, \quad \mathbb{E}[\bar{\Sigma}_{ML}] = \Sigma$$

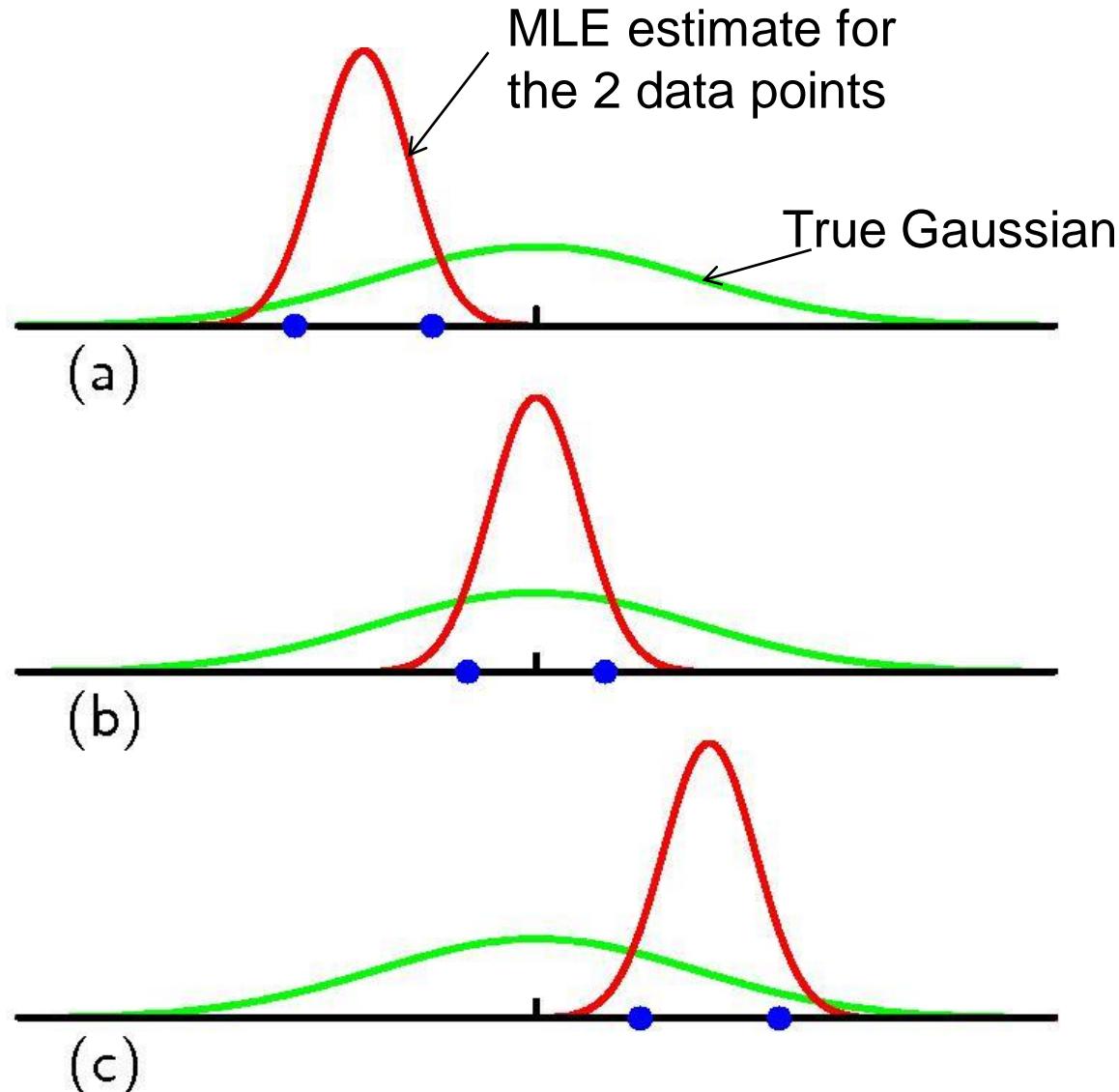
- To prove this, you will need to use that:

$$\mathbb{E}[x_n x_m^T] = \mu \mu^T + \delta_{mn} \Sigma$$



Bias in MLE

- In the schematic from [Bishop's PRML](#), we consider 3 cases each with 2 data points extracted from the true Gaussian.



- The mean of the three distributions predicted via MLE (i.e. averaged over the data) is correct.
- However, the variance is underestimated since it is a variance with respect to the sample mean and NOT the true mean.

Robbins-Monro Algorithm

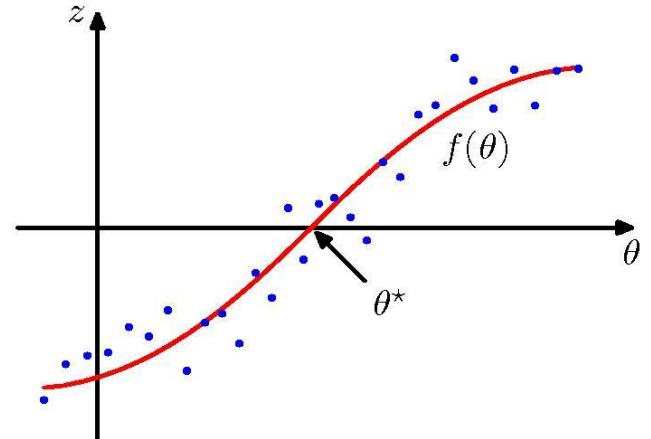
$$f(\theta) = \mathbb{E}(z|\theta) = \int z p(z|\theta) dz$$

- We want to find the root $f(\theta^*) = 0$ in a sequential manner: The Robbins-Monro algorithm proceeds as:

$$\theta^{(N)} = \theta^{(N-1)} - a_{N-1} z(\theta^{(N-1)})$$

- The learning coefficients $\{a_N\}$ should satisfy:

$$\lim_{N \rightarrow \infty} a_N = 0, \sum_{n=1}^{\infty} a_N = \infty, \sum_{n=1}^{\infty} a_N^2 < \infty$$

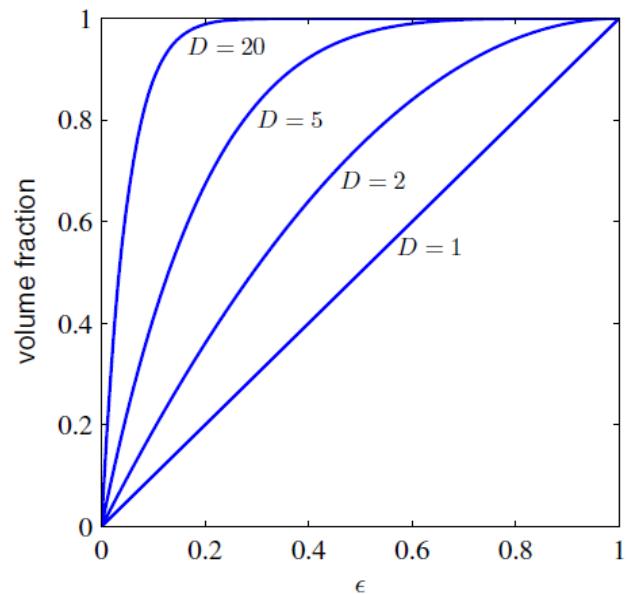


Volume of a Sphere in High Dimensions

- Consider a sphere of radius $r = 1$ in D dimensions. Let us compute the fraction of the volume of the sphere that lies between radius $r = 1 - \varepsilon$ and $r = 1$.

$$\frac{V_D(1) - V_D(1 - \varepsilon)}{V_D(1)} = \frac{K_D 1^D - K_D (1 - \varepsilon)^D}{K_D 1^D} \\ = 1 - (1 - \varepsilon)^D$$

- Note that for large D , this fraction tends to 1 even for small ε
- In high dimensions, the volume of the sphere is concentrated near the surface!



The volume of a sphere in D -dimensions is given as

$$V_D(r) = K_D r^D,$$

K_D = D - dependent constant

Bayes' Theorem and Gaussian Linear Models

- Consider a linear Gaussian model: A Gaussian marginal distribution $p(\mathbf{x})$ and a Gaussian conditional distribution $p(\mathbf{y}|\mathbf{x})$ in which $p(\mathbf{y}|\mathbf{x})$ has a mean that is a linear function of \mathbf{x} , and a covariance which is independent of \mathbf{x} .

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$

$$p(\mathbf{y} | \mathbf{x}) = \mathcal{N}(\mathbf{y} | \mathbf{Ax} + \mathbf{b}, \mathbf{L}^{-1})$$

- We want using Bayes' rule to find $p(\mathbf{y})$ and $p(\mathbf{x}|\mathbf{y})$.
- We start with the joint distribution over $\mathbf{z}=(\mathbf{x},\mathbf{y})$ which is quadratic in the components of \mathbf{z} – so $p(\mathbf{z})$ is a Gaussian.

Posterior Inference: Point Estimates

$$\pi(\theta | x) = \frac{f(x | \theta)\pi(\theta)}{m(x)}$$

Maximum A Posteriori estimate (MAP)

$$\theta^* = \arg \max_{\theta} \log(\pi(\theta | x)) = \arg \max_{\theta} (\log \pi(x | \theta) + \log \pi(\theta))$$

Posterior Mean

$$\hat{\theta} = \mathbb{E}_{p(\theta|x)}[\theta] = \int \theta \pi(\theta | x) d\theta$$

Posterior Quantiles

$$\Pr[\theta > a] = \int_a^{\infty} \pi(\theta | x) d\theta$$



Prediction

Suppose we have observed \mathbf{x} and we want to make a prediction about (future) unknown observables: What is the probability of observing data $\hat{\mathbf{x}}$?
If we already have observed data \mathbf{x} ?

This means finding $g(\hat{\mathbf{x}}|\mathbf{x})$

We have:

$$\begin{aligned} g(\hat{\mathbf{x}}|\mathbf{x}) &= \int g(\hat{\mathbf{x}}, \theta | \mathbf{x}) d\theta = \int \frac{\pi(\hat{\mathbf{x}}, \theta, \mathbf{x})}{m(\mathbf{x})} d\theta = \int \frac{\pi(\hat{\mathbf{x}}, \theta, \mathbf{x})}{\phi(\theta, \mathbf{x})} \frac{\phi(\theta, \mathbf{x})}{m(\mathbf{x})} d\theta = \\ &= \int f(\hat{\mathbf{x}}|\theta, \mathbf{x}) \pi(\theta|\mathbf{x}) d\theta = \int f(\hat{\mathbf{x}}|\theta) \pi(\theta|\mathbf{x}) d\theta \end{aligned}$$

Compare this with the normalizing factor:

$$m(\hat{\mathbf{x}}) = \int f(\hat{\mathbf{x}}|\theta) \pi(\theta) d\theta$$



Marginal Likelihood or Evidence

➤ In the Bayesian parametric model, we define the following:

- The joint distribution of (θ, X)

$$\phi(\theta, x) = \pi(\theta)f(x | \theta)$$

- The marginal distribution of X

$$m(x) = \int \phi(\theta, x)d\theta = \int \pi(\theta)f(x | \theta)d\theta$$

- For a realization $X=x$, $m(x)$ is called **marginal likelihood or evidence**



Normalizing Factor in Bayes' Rule

$$\pi(\theta | x) = \frac{f(x | \theta)\pi(\theta)}{m(x)}$$

- What about $m(x)$ (the marginal distribution on the data x)

$$m(x) = \int \phi(x, \theta) d\theta = \int f(x | \theta)\pi(\theta) d\theta$$

- It is the normalizing *constant of the posterior* $\pi(\theta | x)$
- It is not needed for inference or prediction
- It is essential for model validation (we will discuss this later on)

Predictive Distribution and Approximations

- Given the prior $\pi(\theta)$ and the likelihood $\ell(\theta|x) = f(x|\theta)$, Bayes' formula yields:^a

$$\pi(\theta|x) = \frac{f(x|\theta)\pi(\theta)}{\int f(x|\theta)\pi(\theta)d\theta}$$

- This represents all the information on θ that can be extracted from x .
- Note the integral at the denominator of the Bayes' rule.
- The predictive distribution of Y when $Y \sim g(y|\theta, x)$ is (we will revisit this in a forthcoming lecture)

$$g(y|x) = \int g(y, \theta|x)d\theta = \underbrace{\int g(y|\theta, x)}_{Likelihood} \underbrace{\pi(\theta|x)}_{Posterior} d\theta$$

- This is to distinguish from prediction using $\hat{\theta} = \theta^{MLE}$ or $\hat{\theta} = \theta^{MAP}$: $g(y|\hat{\theta}, x)$

^a [Probability, Conditional Probability and Bayes Formula](#), B. Vidakovic



Exponential Family

- The exponential family of distributions over x , given parameters η , is defined to be the set of distributions of the form

$$p(x | \eta) = h(x)g(\eta)\exp\{\eta^T u(x)\} \text{ or}$$

$$p(x | \eta) = h(x)\exp\{\eta^T u(x) - A(\eta)\}, \text{ where } A(\eta) = -\log g(\eta)$$

x is scalar/vector, discrete/continuous. **η are the natural parameters and $u(x)$ is referred to as a sufficient statistic.**

- $g(\eta)$ ensures that the distribution is normalized and satisfies

$$g(\eta) \int h(x) \exp\{\eta^T u(x)\} dx = 1$$

- The normalization factor Z and the log of it A are defined as:

$$Z(\eta) = \frac{1}{g(\eta)}, A(\eta) = \ln Z(\eta) = -\ln g(\eta) = \ln \int h(x) \exp\{\eta^T u(x)\} dx$$

$$p(x | \eta) = h(x) \exp\{\eta^T u(x)\} / Z(\eta)$$

- The space of η for which $\int h(x) \exp\{\eta^T u(x)\} dx < \infty$ is the **natural parameter space**.



Conjugate Priors

- For any member of the exponential family,

$$p(x | \theta) = h(x)g(\eta(\theta))\exp\{\eta^T(\theta)u(x)\}$$

there exists a conjugate prior that can be written in the form

$$p(\theta | \nu_0, \tau_0) \propto g(\eta(\theta))^{\nu_0} \exp\{\eta^T(\theta)\tau_0\} = \exp\{\nu_0 \eta^T(\theta)\bar{\tau}_0 - A(\eta(\theta))\nu_0\}, \text{ where: } \tau_0 \equiv \nu_0 \bar{\tau}_0$$

- In normalized form, we write:

$$p(\theta | \nu_0, \tau_0) = \frac{1}{Z(\nu_0, \tau_0)} g(\eta(\theta))^{\nu_0} \exp\{\eta^T(\theta)\tau_0\} = \frac{1}{Z(\nu_0, \tau_0)} \exp\{\nu_0 \eta^T(\theta)\bar{\tau}_0 - A(\eta(\theta))\nu_0\}$$

$$\text{where: } Z(\nu_0, \tau_0) = \int \exp\{\nu_0 \eta^T(\theta)\bar{\tau}_0 - A(\eta(\theta))\nu_0\} d\theta$$

Standard Exponential Families

$f(x \theta)$	$\pi(\theta)$	$\pi(\theta x)$
Normal $N(\theta, \sigma^2)$	Normal $\mathcal{N}(\mu, \tau^2)$	$\mathcal{N}(\rho(\sigma^2\mu + \tau^2x), \rho\sigma^2\tau^2)$ $\rho^{-1} = \sigma^2 + \tau^2$
Poisson $\mathcal{P}(\theta)$	Gamma $\mathcal{P}(\theta)\mathcal{G}(\alpha, \beta)$	$\mathcal{G}(\alpha + x, \beta + 1)$
Gamma $\mathcal{G}(v, \theta)$	Gamma $\mathcal{G}(\alpha, \beta)$	$\mathcal{G}(\alpha + v, \beta + x)$
Binomial $\mathcal{B}(n, \theta)$	Beta $\mathcal{B}\text{e}(\alpha, \beta)$	$\mathcal{B}\text{e}(\alpha + x, \beta + n - x)$
Negative Binomial $\mathcal{N}\text{eg}(m, \theta)$	Beta $\mathcal{B}\text{e}(\alpha, \beta)$	$\mathcal{B}\text{e}(\alpha + m, \beta + x)$
Multinomial $\mathcal{M}_k(\theta_1, \dots, \theta_k)$	Dirichlet $\mathcal{D}(\alpha_1, \dots, \alpha_k)$	$\mathcal{D}(\alpha_1 + x_1, \dots, \alpha_k + x_k)$
Normal $\mathcal{N}(\mu, 1/\theta)$	Gamma $\mathcal{G}\text{a}(\alpha, \beta)$	$\mathcal{G}(\alpha + 0.5, \beta + (\mu - x)^2/2)$

Generalizing to Multivariate Gaussians

- Suppose you have $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \sim (\text{i.i.d}) \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- We do not know $\boldsymbol{\mu}$ or $\boldsymbol{\Sigma}$
- When both sets of parameters are unknown, a conjugate family of priors is one in which

$$\boldsymbol{\Sigma} \sim \mathcal{IW}(\nu, \boldsymbol{\Lambda}^{-1})$$

and

$$\boldsymbol{\mu} | \boldsymbol{\Sigma} \sim \mathcal{N}(\boldsymbol{\eta}, \boldsymbol{\Sigma} / \kappa)$$

- The Wishart distribution is a multivariate analog of the Gamma distribution. If matrix \mathbf{U} has the Wishart distribution, then \mathbf{U}^{-1} has the inverse-Wishart distribution.
- The quantity ν is a positive scalar, while $\boldsymbol{\Lambda}$ is a positive definite matrix. They play roles analogous to those played by α and β , respectively, in the Gamma distribution.
- The other parameters of the prior are the mean vector $\boldsymbol{\eta}$ and κ , the latter of which represents the ``a priori number of observations''.



Wishart Distribution

Wishart	$W \sim \text{Wishart}_\nu(S)$ $p(W) = \text{Wishart}_\nu(W S)$ (implicit dimension $k \times k$)	degrees of freedom ν symmetric, pos. definite $k \times k$ scale matrix S
---------	--	---

$$p(W) = \left(2^{\nu k/2} \pi^{k(k-1)/4} \prod_{i=1}^k \Gamma\left(\frac{\nu+1-i}{2}\right) \right)^{-1} \times |S|^{-\nu/2} |W|^{(\nu-k-1)/2} \times \exp\left(-\frac{1}{2}\text{tr}(S^{-1}W)\right), W \text{ pos. definite}$$

Gamma	$\theta \sim \text{Gamma}(\alpha, \beta)$ $p(\theta) = \text{Gamma}(\theta \alpha, \beta)$	shape $\alpha > 0$ inverse scale $\beta > 0$
-------	---	---

$$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta}, \quad \theta > 0$$
$$\begin{aligned} \text{E}(\theta) &= \frac{\alpha}{\beta} \\ \text{var}(\theta) &= \frac{\alpha}{\beta^2} \\ \text{mode}(\theta) &= \frac{\alpha-1}{\beta}, \text{ for } \alpha \geq 1 \end{aligned}$$

Bayesian Data Analysis, A. Gelman, J. Carlin, H. Stern and D. Rubin, 2004



Jeffrey's Noninformative Priors

- Jeffrey's proposes a more intrinsic approach which avoids the need to take the invariance structure into account.
- Given a likelihood $f(x | \theta)$, Jeffrey's noninformative prior distributions are based on **Fisher information**, given by

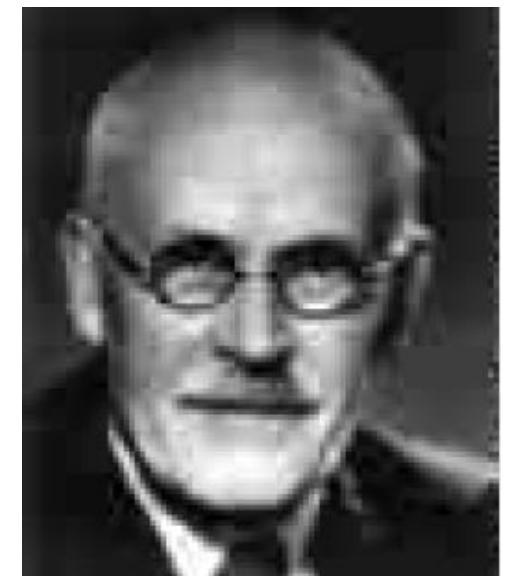
$$I(\theta) = \mathbb{E}_{X|\theta} \left(\frac{\partial \log f(X | \theta)}{\partial \theta} \frac{\partial \log f(X | \theta)^T}{\partial \theta} \right) = -\mathbb{E}_{X|\theta} \left(\frac{\partial^2 \log f(X | \theta)}{\partial \theta^2} \right)$$

the corresponding prior distribution is

$$\pi(\theta) \propto |I(\theta)|^{-1/2}$$

Determinant of I

Sir Harold Jeffreys
(1891–1989)



Empirical Bayes - Evidence Approximation

- In hierarchical Bayesian models, we need to compute the posterior on multiple levels of latent variables. For example, in a two-level model, we need to compute

$$p(\boldsymbol{\eta}, \boldsymbol{\theta} | \mathcal{D}) \propto p(\mathcal{D} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \boldsymbol{\eta}) p(\boldsymbol{\eta})$$

- In some cases, we can analytically marginalize out $\boldsymbol{\theta}$; this leaves us with the simpler problem of just computing $p(\boldsymbol{\eta} | \mathcal{D})$.
- As a computational shortcut, we can *approximate the posterior on the hyper-parameters with a point-estimate*,

$$p(\boldsymbol{\eta} | \mathcal{D}) \approx \delta_{\bar{\boldsymbol{\eta}}}(\boldsymbol{\eta}), \bar{\boldsymbol{\eta}} = \operatorname{argmax} p(\boldsymbol{\eta} | \mathcal{D}) = \operatorname{argmax} \left[\int p(\mathcal{D} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \boldsymbol{\eta}) d\boldsymbol{\theta} \right]$$

- Since $\boldsymbol{\eta}$ is typically much smaller than $\boldsymbol{\theta}$ in dimensionality, it is less prone to overfitting, so we can safely use a uniform prior on $\boldsymbol{\eta}$.
- The quantity inside the brackets is *the marginal or integrated likelihood, often called the evidence*. This overall approach is called *empirical Bayes (EB) or type-II maximum likelihood*. In machine learning, it is sometimes called the *evidence procedure*.



Empirical Bayes

- Empirical Bayes violates the principle that the prior should be chosen independently of the data.
- We can just view it as a cheap approximation to inference in a hierarchical Bayesian model, just as we viewed MAP estimation as an approximation to inference in the one level model $\theta \rightarrow \mathcal{D}$.
- We can construct a hierarchy in which the more integrals one performs, the “more Bayesian” one becomes:

Method	Definition
Maximum likelihood	$\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathcal{D} \theta)$
MAP estimation	$\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathcal{D} \theta)p(\theta \eta)$
ML-II (Empirical Bayes)	$\hat{\eta} = \operatorname{argmax}_{\eta} \int p(\mathcal{D} \theta)p(\theta \eta)d\theta = \operatorname{argmax}_{\eta} p(\mathcal{D} \eta)$
MAP-II	$\hat{\eta} = \operatorname{argmax}_{\eta} \int p(\mathcal{D} \theta)p(\theta \eta)p(\eta)d\theta = \operatorname{argmax}_{\eta} p(\mathcal{D} \eta)p(\eta)$
Full Bayes	$p(\theta, \eta \mathcal{D}) \propto p(\mathcal{D} \theta)p(\theta \eta)p(\eta)$

Bayesian Model Validation

From data we can learn the parameters for each model
and then the model itself

$$x \Rightarrow \pi_i(\theta_i | x, M_i) = \frac{f_i(x | \theta_i, M_i) \pi_i(\theta_i | M_i)}{\pi_i(x | M_i)} \Rightarrow \pi_i(M_i | x) = \frac{\pi_i(x | M_i) \pi_i(M_i)}{\pi(x)}$$

Noting that

$$\pi_i(x | M_i) = \int f_i(x | \theta_i, M_i) \pi_i(\theta_i | M_i) d\theta_i$$

we can find the best model that represents the data by computing:

$$\frac{\pi(M_1 | x)}{\pi(M_2 | x)} = \frac{\pi(x | M_1) \pi(M_1)}{\pi(x | M_2) \pi(M_2)} = \underbrace{\frac{\int f_1(x | \theta_1, M_1) \pi_1(\theta_1 | M_1) d\theta_1}{\int f_2(x | \theta_2, M_2) \pi_2(\theta_2 | M_2) d\theta_2}}_{\text{Ratio of Bayes' factors}} \underbrace{\frac{\pi(M_1)}{\pi(M_2)}}_{\text{Ratio of Priors}}$$



Bayesian Occam's Razor

- To further understand the Bayesian Occam's razor effect is to note that probabilities must sum to one (sum over all possible data sets)

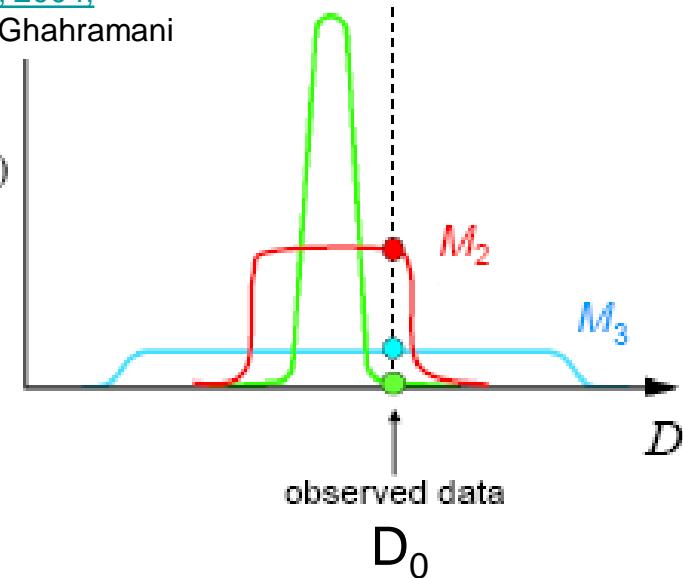
$$\sum_{\mathcal{D}'} p(\mathcal{D}' | m) = 1$$

[Bayesian Methods for
Machine Learning, ICML
Tutorial, 2004,](#)
Zoubin Ghahramani

- Model 1 is too simple and assigns low probability to D_0 .

$$p(D = d | M)$$

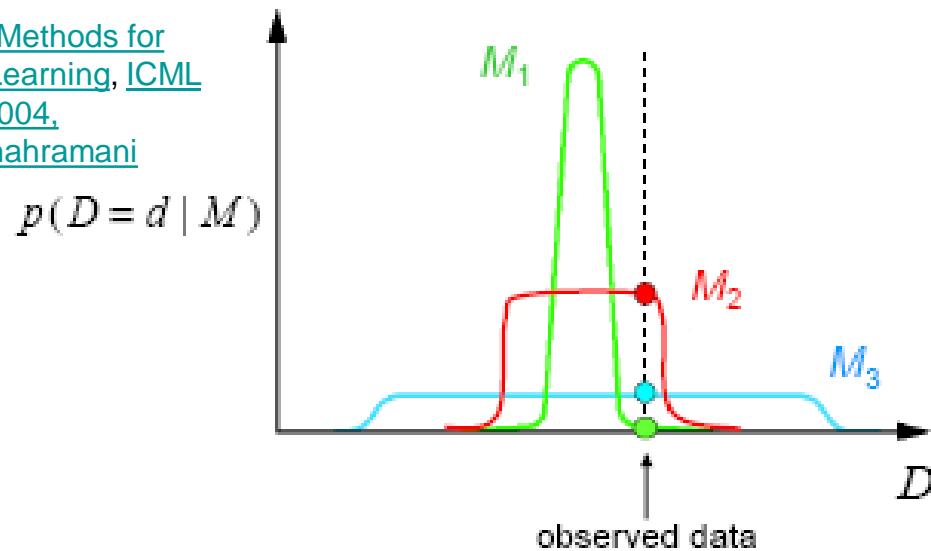
- Model 3 also assigns D_0 relatively low probability, because it can predict many data sets, and hence it spreads its probability quite widely & thinly.



- Model 2 is "just right": it predicts the observed data with a reasonable degree of confidence, but does not predict too many other things. Hence model 2 is the most probable model.

Bayesian Occam's Razor

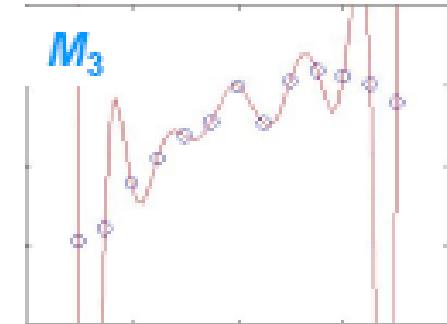
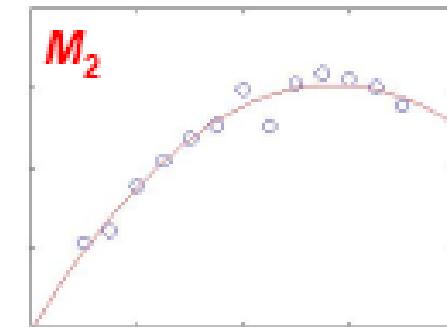
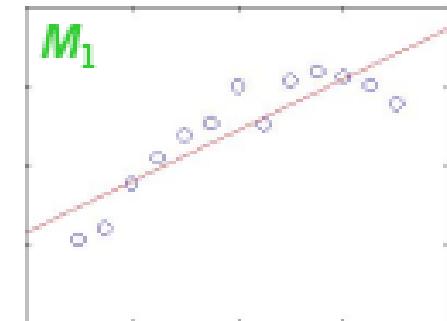
Bayesian Methods for
Machine Learning, ICML
Tutorial, 2004,
Zoubin Ghahramani



M_1 : the too simple model is unlikely to generate this data

M_3 : the too complex model explains poorly a lots of data sets and it is a little better but still unlikely to have generated our data

M_2 : the just right model has the highest marginal likelihood



Laplace Approximation

- The Laplace approximation allows a Gaussian approximation of the parameter posterior about the maximum a posteriori (MAP) parameter estimate.
- Consider a data set \mathcal{D} and M models $\mathcal{M}_i, i=1,\dots,M$ with corresponding parameters $\theta_i, i=1,\dots,M$. We compare models using the posteriors:

$$p(\mathcal{M} | \mathcal{D}) \propto p(\mathcal{M}) p(\mathcal{D} | \mathcal{M})$$

- For large sets of data \mathcal{D} (relative to the model parameters), the parameter posterior is approximately Gaussian around the MAP estimate θ_m^{MAP} (can also use 2nd order Taylor expansion of the log-posterior):

$$p(\theta_m | \mathcal{D}, \mathcal{M}_m) \approx (2\pi)^{-d/2} |A|^{1/2} \exp\left(-\frac{1}{2} (\theta_m - \theta_m^{MAP})^T A (\theta_m - \theta_m^{MAP})\right),$$
$$A_{ij} = -\left. \frac{\partial^2 \log P(\theta_m | \mathcal{D}, \mathcal{M}_m)}{\partial \theta_{mi} \partial \theta_{mj}} \right|_{\theta_m^{MAP}}$$



Sampling From an Arbitrary Distribution

- Consider an arbitrary probability density $\pi(\mathbf{x})$
- Monte Carlo approximation is given by

$$\hat{\pi}_N(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \delta_{X^{(i)}}(\mathbf{x}), \text{ where } X^{(i)} \stackrel{i.i.d.}{\sim} \pi$$

- For any function $f : \mathcal{X} \rightarrow \mathbb{R}$

$$\mathbb{E}_{\hat{\pi}_N}(f) = \frac{1}{N} \sum_{i=1}^N f(X^{(i)}) \cong \mathbb{E}_{\pi}(f)$$

But how do
we sample from
an arbitrary
distribution?



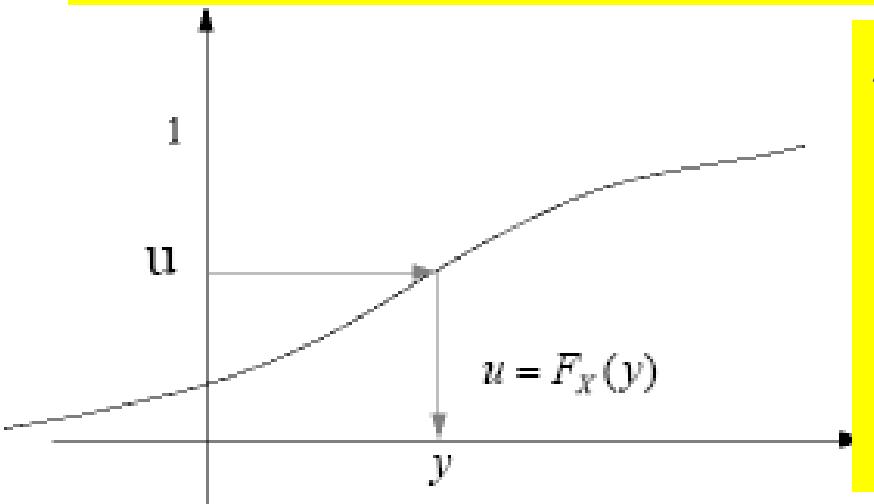
or more precisely:

$$\mathbb{E}_X[\mathbb{E}_{\hat{\pi}_N}(f)] = \mathbb{E}_{\pi}(f) \text{ and } \text{Var}_X(\mathbb{E}_{\hat{\pi}_N}(f)) = \frac{\text{Var}_{\pi}(f)}{N}$$



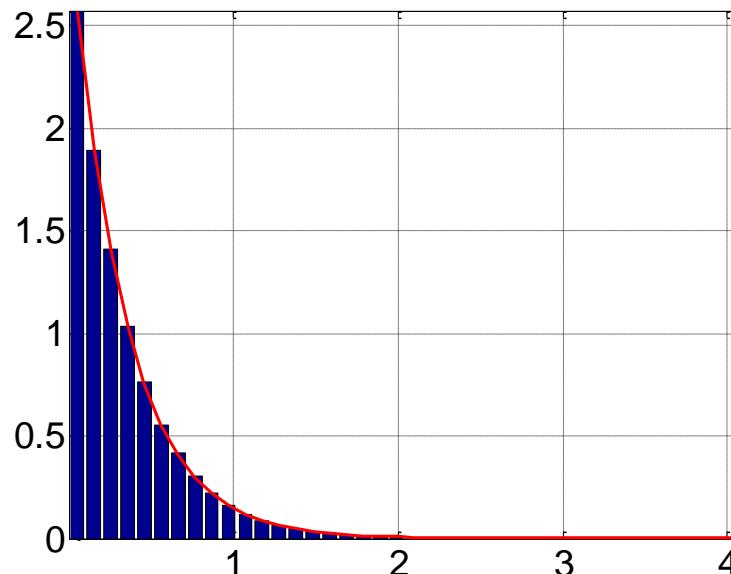
Inverse Method: Exponential Distribution

Exponential(λ): $\pi(x) = \lambda e^{-\lambda x}$, $F_X(y) = 1 - e^{-\lambda y}$, $x = -\ln(1-u)/\lambda$



As $F(y)$ is area under $\pi(y)$, $y = F^{-1}(u)$ prescribes that

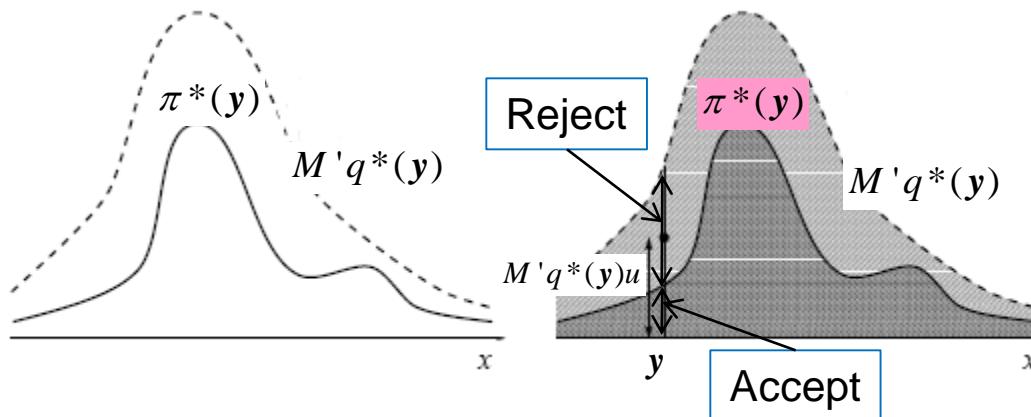
- Choose $u=(0,1]$, then find value y that has that fraction u of area to the left of y , or $u=F(y)$
- Return that value of $x=y$.



[MatLab Implementation](#)



Rejection Sampling



- Set $i=1$
- Repeat until $i=N$
 - Sample $y \sim q(y)$ and $u \sim \mathcal{U}_{(0,1)}$
 - If $u < \frac{\pi^*(y)}{M'q^*(y)}$ then accept (set $x^{(i)}=y$) and increment the counter i
 - ✓ Otherwise, reject

The distribution $\pi(y)$ needs to be known only up to a normalizing constant:

$$\pi(y) = \frac{\pi^*(y)}{Z}, Z = \int \pi^*(y) dy$$

Importance Sampling

- Computing $\mathbb{E}_\pi\{f(X)\}$ is transformed to computing

$$\mathbb{E}_q \left\{ \frac{f(X)\pi(X)}{\underbrace{q(X)}_{f_q(X)}} \right\}$$

- We can employ MC for this task:

- Generate samples $X^{(i)} \sim q(x)$ (i.i.d.) and evaluate:

$$\mathbb{E}_{\hat{q}_N}(f_q(x)) = \frac{1}{N} \sum_{i=1}^N f_q(X^{(i)}) = \frac{1}{N} \sum_{i=1}^N \frac{\pi(X^{(i)})}{q(X^{(i)})} f(X^{(i)}) = \frac{1}{N} \sum_{i=1}^N w(X^{(i)}) f(X^{(i)})$$

where the importance weight is $w(X^{(i)}) \equiv \frac{\pi(X^{(i)})}{q(X^{(i)})}$.

Asymptotic Variance of the Normalized IS

- It follows that

$$\begin{aligned} \text{Var}\left(\mathbb{E}_{\hat{\pi}_N}(f(X))\right) &\simeq \sigma_A^2 \left(\frac{\partial g}{\partial A}(\mu)\right)^2 + \sigma_B^2 \left(\frac{\partial g}{\partial B}(\mu)\right)^2 + 2 \frac{\partial g}{\partial A}(\mu) \frac{\partial g}{\partial B}(\mu) \sigma_{A,B} = \\ &= \frac{\sigma_A^2}{\mu_B^2} + \frac{\sigma_B^2 \mu_A^2}{\mu_B^4} - 2 \frac{\sigma_{A,B} \mu_A}{\mu_B^3} \end{aligned}$$

- Straight forward substitution of the expressions for $\sigma_A^2, \sigma_B^2, \mu_A, \mu_B$, and using $\mathbb{E}_q(\mathbb{E}_{\hat{\pi}_N}(f(X))) \simeq \mathbb{E}_\pi(f(X))$ leads to the following CLT asymptotic result:

$$\sqrt{N} \left(\mathbb{E}_{\hat{\pi}_N}(f(X)) - \mathbb{E}_\pi(f(X)) \right) \sim \mathcal{N}(0, \sigma_{IS}^2(f))$$

where

$$\sigma_{IS}^2(f) = \int \frac{\pi^2(x)}{q(x)} (f(x) - \mathbb{E}_\pi(f))^2 dx = \text{Var}_q(w(X)(f(X) - \mathbb{E}_\pi(f)))$$

Importance Sampling in High-Dimensions

- Consider the following target distribution

$$\pi(\mathbf{x}) = \mathcal{N}(\mathbf{0}, \mathbf{I}) = \frac{1}{(2\pi)^{d/2}} e^{-\frac{\sum_{i=1}^d x_i^2}{2}}$$

- We take the following reasonable importance sampling distribution ($\sigma > 1$)

$$q_\sigma(\mathbf{x}) = \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) = \frac{1}{(2\pi\sigma^2)^{d/2}} e^{-\frac{\sum_{i=1}^d x_i^2}{2\sigma^2}}$$

- Note that:

$$w_\sigma(x) = \frac{\pi(x)}{q_\sigma(x)} = \sigma^d e^{-\frac{1}{2}\sum_{i=1}^d x_i^2(1-\frac{1}{\sigma^2})} \leq \sigma^d \quad \forall x$$

Importance Sampling in High-Dimensions

$$q_\sigma(x) = \mathcal{N}(\theta, \sigma^2 I) = \frac{1}{(2\pi\sigma^2)^{d/2}} e^{-\frac{\sum_{i=1}^d x_i^2}{2\sigma^2}}$$
$$w_\sigma(x) = \frac{\pi(x)}{q_\sigma(x)} = \sigma^d e^{-\frac{1}{2}\sum_{i=1}^d x_i^2(1-\frac{1}{\sigma^2})} \leq \sigma^d \quad \forall x$$

□ We now note that:

$$\begin{aligned} Var_{q_\sigma}\left(\frac{\pi(x)}{q_\sigma(x)}\right) &= \int q_\sigma \frac{\pi^2(x)}{q_\sigma^2(x)} dx - \left(\int \pi(x) dx\right)^2 = \\ &= \int \frac{1}{(2\pi)^{d/2}} \sigma^d \exp\left\{-\frac{\sum_{i=1}^d x_i^2}{2}\left(2 - \frac{1}{\sigma^2}\right)\right\} dx - 1 = \sigma^d \left(\frac{\sigma^2}{2\sigma^2 - 1}\right)^{d/2} - 1 = \underbrace{\left(\frac{\sigma^4}{2\sigma^2 - 1}\right)^{d/2}}_{>1} - 1 \end{aligned}$$

□ It is easy to see that: $\sigma^4 > 2\sigma^2 - 1 \Leftrightarrow (\sigma^2 - 1)^2 > 0$. Therefore:

$$Var_{q_\sigma}\left(\frac{\pi(x)}{q_\sigma(x)}\right) \rightarrow \infty \text{ as } d \rightarrow \infty$$

□ The variance of the weights increases exponentially fast with dimensionality. This is despite the good choice of $q(x)$.



Introduction to Markov Chain Monte Carlo

- **Markov chain:** A sequence of random variables $\{X_n, n \in \mathbb{N}\}$ defined on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ such that for any $A \in \mathcal{B}(\mathcal{X})$ the following probability condition is satisfied:

$$\mathbb{P}(X_n \in A | X_0, \dots, X_{n-1}) = \mathbb{P}(X_n \in A | X_{n-1})$$

and we write:

$$\text{Transition Kernel : } P(x, A) = \mathbb{P}(X_n \in A | X_{n-1})$$

- **Markov Chain Monte Carlo (MCMC):** Given a target distribution π , we need to design a transition kernel P such that asymptotically

$$\frac{1}{N} \sum_{n=1}^N f(X_n) \xrightarrow{N \rightarrow \infty} \int f(x) \pi(x) dx \text{ and / or } X_n \sim \pi$$

- It is easy to simulate the Markov Chain even if π is complex.



Gibbs Sampler

- If $\theta = (\theta_1, \theta_2, \dots, \theta_p)$ where $p > 2$, the Gibbs sampler still applies.
- Initialization:
 - Select deterministically or randomly $\theta^{(0)} = (\theta_1^{(0)}, \theta_2^{(0)}, \dots, \theta_p^{(0)})$
- Iteration i , $i \geq 1$
 - For $k=1:p$
 - Sample $\theta_k^{(i)} \sim \pi(\theta_k | \theta_{-k}^{(i)})$
 - where $\theta_{-k}^{(i)} = (\theta_1^{(i)}, \dots, \theta_{k-1}^{(i)}, \theta_{k+1}^{(i)}, \dots, \theta_p^{(i)})$



Metropolis -Hastings Algorithm

- Let $\pi(x)$ the target and $q(y | x)$ any (symmetric or not) distribution such $q(y | x) = q(x | y)$. Given state x_n at step n

- Draw a proposal y from $q(y | x_n)$

- Calculate acceptance ratio:

$$\alpha(x_n, y) = \min \left\{ 1, \frac{\pi(y)}{\pi(x_n)} \frac{q(x_n | y)}{q(y | x_n)} \right\}$$

- Set

$$x_{n+1} = \begin{cases} y & \text{with probability } \alpha(x_n, y) \\ x_n & \text{with probability } 1 - \alpha(x_n, y) \end{cases}$$

W. Hastings, [Monte Carlo Sampling Methods using Markov Chains and their Applications](#), Biometrika, Vol. 57(1), pp. 97-109 (1970).



Composition of MH Kernels

- Assume we use a composition of these kernels, then the resulting algorithm proceeds as follows at iteration i.

MH Step to Update Component 1

- *Sample $\theta_1^* \sim q_1\left(\left(\theta_1^{(i-1)}, \theta_2^{(i-1)}\right), \cdot\right)$ and compute*

$$\alpha_1\left(\left(\theta_1^{(i-1)}, \theta_2^{(i-1)}\right), \left(\theta_1^*, \theta_2^{(i-1)}\right)\right) = \min\left(1, \frac{\pi\left(\theta_1^* \mid \theta_2^{(i-1)}\right) q_1\left(\left(\theta_1^*, \theta_2^{(i-1)}\right), \theta_1^{(i-1)}\right)}{\pi\left(\theta_1^{(i-1)} \mid \theta_2^{(i-1)}\right) q_1\left(\left(\theta_1^{(i-1)}, \theta_2^{(i-1)}\right), \theta_1^*\right)}\right)$$

- *With probability $\alpha_1\left(\left(\theta_1^{(i-1)}, \theta_2^{(i-1)}\right), \left(\theta_1^*, \theta_2^{(i-1)}\right)\right)$, set $\theta_1^{(i)} = \theta_1^*$;*

otherwise set $\theta_1^{(i)} = \theta_1^{(i-1)}$

Composition of MH Kernels

- Assume we use a composition of these kernels, then the resulting algorithm proceeds as follows at iteration i.

MH Step to Update Component 2

- *Sample $\theta_2^* \sim q_2\left(\left(\theta_1^{(i)}, \theta_2^{(i-1)}\right), \cdot\right)$ and compute*

$$\alpha_2\left(\left(\theta_1^{(i)}, \theta_2^{(i-1)}\right), \left(\theta_1^{(i)}, \theta_2^*\right)\right) = \min\left(1, \frac{\pi\left(\theta_2^* \mid \theta_1^{(i)}\right) q_2\left(\left(\theta_1^{(i)}, \theta_2^*\right), \theta_2^{(i-1)}\right)}{\pi\left(\theta_2^{(i-1)} \mid \theta_2^{(i)}\right) q_1\left(\left(\theta_1^{(i)}, \theta_2^{(i-1)}\right), \theta_2^{(i)}\right)}\right)$$

- *With probability $\alpha_2\left(\left(\theta_1^{(i)}, \theta_2^{(i-1)}\right), \left(\theta_1^{(i)}, \theta_2^*\right)\right)$, set $\theta_2^{(i)} = \theta_2^*$;*

otherwise set $\theta_2^{(i)} = \theta_2^{(i-1)}$



Sequential Importance Sampling

At step n=1:

- Sample $X_1^{(i)} \sim q(x_1 | y_1), i = 1, \dots, N$ and then approximate:

$$\hat{p}_N(x_1 | y_1) = \sum_{i=1}^N W_1^{(i)} \delta_{X_1^{(i)}}(x_1), \quad W_1^{(i)}(X_1^{(i)}, y_1) \propto \frac{\mu(X_1^{(i)}) g(y_1, X_1^{(i)})}{q(X_1^{(i)} | y_1)}$$

At step $n \geq 2$:

- Sample $X_n^{(i)} \sim q(x_n | y_n, X_{n-1}^{(i)}), n = 1, \dots, N$ and compute:

$$\hat{p}_N(\mathbf{x}_{1:n} | \mathbf{y}_{1:n}) = \sum_{i=1}^N W_n^{(i)} \delta_{X_{1:n}^{(i)}}(\mathbf{x}_{1:n}),$$

$$W_n^{(i)} \propto w(\mathbf{X}_{1:n}^{(i)}, \mathbf{y}_{1:n}) = w(\mathbf{X}_{1:n-1}^{(i)}, \mathbf{y}_{1:n-1}) \frac{f(X_n^{(i)} | X_{n-1}^{(i)}) g(y_n | X_n^{(i)})}{q(X_n^{(i)} | y_n, X_{n-1}^{(i)})}$$

- The algorithm has computational complexity $\mathcal{O}(N)$ independent of n.



Sequential Importance Sampling

- Note that the complexity of the algorithm does not increase with n .
- The algorithm is fully parallelizable.
- Also note that if our interest is on computing the marginal posterior,

$\hat{p}_N(x_n | \mathbf{y}_{1:n})$ (posterior filtered density), then we only need to store $X_{n-1:n}^{(i)}$ rather than all the $X_{1:n}^{(i)}$ paths

$$\hat{p}_N(x_n | \mathbf{y}_{1:n}) = \sum_{i=1}^N W_n^{(i)} \delta_{X_n^{(i)}}(x_n),$$
$$W_n^{(i)} \propto w\left(\mathbf{X}_{1:n}^{(i)}, \mathbf{y}_{1:n}\right) = w\left(\mathbf{X}_{1:n-1}^{(i)}, \mathbf{y}_{1:n-1}\right) \frac{f\left(X_n^{(i)} | X_{n-1}^{(i)}\right) g\left(y_n | X_n^{(i)}\right)}{q\left(X_n^{(i)} | y_n, X_{n-1}^{(i)}\right)}$$

- One can show that this approaches the true posterior as $N \rightarrow \infty$.
 - Crisan, D., P. D. Moral, and T. Lyons (1999). [Discrete filtering using branching and interacting particle systems](#). *Markov Processes and Related Fields* 5(3), 293–318.

Conditional Dynamic Linear Model

- A conditional dynamic linear model (CDLM) can be generally defined as

$$Z_n = A(X_n)Z_{n-1} + B(X_n)V_n \quad , \quad Z_1 \sim \mathcal{N}(m_1, \Sigma_1), \quad V_n \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \Sigma_v)$$

$$Y_n = C(X_n)Z_n + D(X_n)W_n \quad , \quad W_n \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \Sigma_w)$$

where all coefficient matrices are functions of the indicator random variable X_n . Both $\{X_n, Z_n\}, n \geq 1$, are unobserved.

- The X_n , which can be either continuous or discrete, is a latent indicator process with certain probabilistic structure. For example:

$$X_1 \sim \mu, \quad X_n | X_{n-1} = x \sim f(\cdot | x)$$

- Here Z_n is an unobserved Markov process and Y_n is the corresponding observation.
- An important feature of the CDLM is that given the trajectory of the indicator variable X_n , the system is Gaussian and linear.

- R.Chen,J.S.Liu, [Mixture Kalman Filters](#), [Journal of the Royal Statistical Society](#), 2000,62B:p493



Birth/Death Moves

- Assume a distribution defined on $\{1\} \times \mathbb{R} \cup \{2\} \times \mathbb{R} \times \mathbb{R}$
- We want to propose moves going from $(1, \theta)$ to $(2, \theta_1, \theta_2)$.
- One can propose

$$u \sim g \in \mathbb{R},$$
$$(\theta_1, \theta_2) = h(\theta, u) = (\theta, u)$$

- Its inverse is given by

$$(\theta, u) = h'(\theta_1, \theta_2) = (\theta_1, \theta_2)$$

- The acceptance probability for this **birth move** is the:

$$\min \left\{ 1, \frac{\pi(2, \theta_1, \theta_2)}{\pi(1, \theta)} \frac{1}{g(u)} \left| \frac{\partial(\theta_1, \theta_2)}{\partial(\theta, u)} \right| \right\} = \min \left\{ 1, \frac{\pi(2, \theta_1, \theta_2)}{\pi(1, \theta_1)} \frac{1}{g(\theta_2)} \right\}$$

Latent Variable Viewpoint of EM

- Initialization: Choose initial set of parameters θ^{old}
- E-step: use current parameters θ^{old} to compute $p(\mathbf{Z}|\mathbf{X}, \theta^{old})$ and then the expected complete-data log-likelihood for general θ

$$Q(\theta, \theta^{old}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{old}) \ln p(\mathbf{X}, \mathbf{Z}|\theta)$$

- M-step: Determine θ^{new} by maximizing $Q(\theta, \theta^{old})$

$$\theta^{new} = \arg \max_{\theta} Q(\theta, \theta^{old})$$

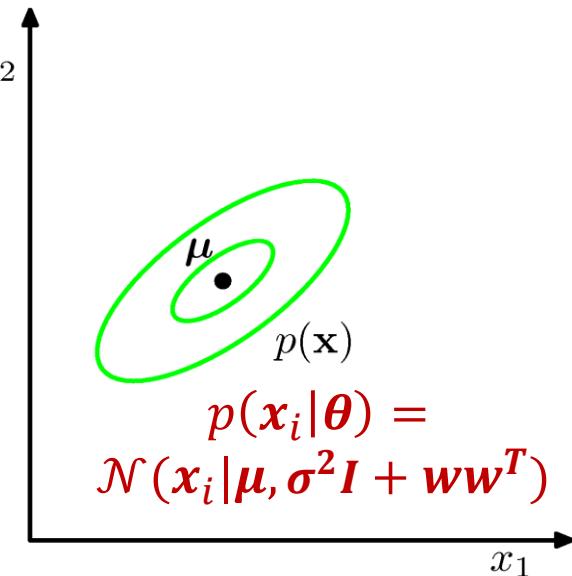
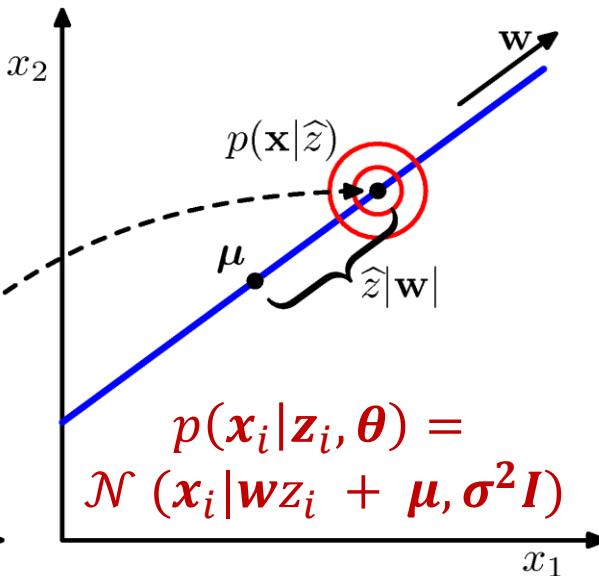
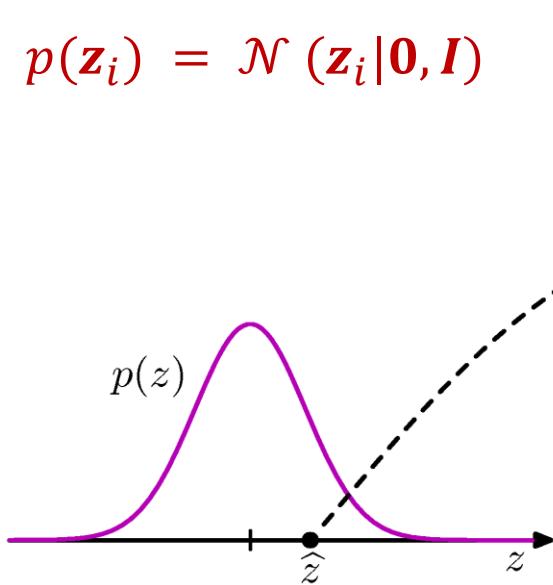
- Check Convergence: stop, or $\theta^{old} \leftarrow \theta^{new}$ and return to E-Step.
- Note: in the definition of $Q(\theta, \theta^{old})$ the log acts directly on $p(\mathbf{X}, \mathbf{Z}|\theta)$ *This step is now computationally tractable.*

Generative Point of View

- PPCA is illustrated for $D = 2, M = 1$ (latent dimension).

$$p(\mathbf{x}_i|\boldsymbol{\theta}) = \int \mathcal{N}(\mathbf{x}_i|W\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}) \mathcal{N}(\mathbf{z}_i|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) d\mathbf{z}_i = \mathcal{N}(\mathbf{x}_i|W\boldsymbol{\mu}_0 + \boldsymbol{\mu}, \boldsymbol{\Psi} + W\boldsymbol{\Sigma}_0 W^T)$$

$$p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i|\mathbf{0}, \mathbf{I})$$



$$p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$$

$$p(\mathbf{x}_i|\mathbf{z}_i, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i|W\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi})$$

$$p(\mathbf{x}_i|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}, \boldsymbol{\Psi} + W\boldsymbol{\Sigma}_0 W^T)$$

- Based on the obtained marginal, we see that without sacrificing generality, we can simplify as:

$$p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i|\mathbf{0}, \mathbf{I})$$

$$p(\mathbf{x}_i|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}, \boldsymbol{\Psi} + W\boldsymbol{\Sigma}_0 W^T)$$